

TECHNICAL MEMORANDUM

STATISTICAL ANALYSIS OF LEAD SAMPLES COLLECTED FROM PROPERTIES LOCATED WITHIN A QUARTER MILE OF THE HERCULANEUM LEAD SMELTER BOUNDARY

Tetra Tech EM Inc. (Tetra Tech) was tasked by the U.S. Environmental Protection Agency (EPA) Region 7 Enforcement Fund Lead Removal program to conduct a trend analysis of soil lead concentrations at selected locations within Herculaneum, Missouri (City). Specifically, the Tetra Tech Superfund Technical Assessment and Response Team (START) 2 was requested to review and analyze data that will enable EPA to determine if soil lead concentrations are increasing and the potential rate of increase at locations within 0.25 mile of the lead smelter boundary. The assessment was conducted under the authority of the Comprehensive Environmental Response, Compensation, and Liability Act of 1980 (CERCLA) and the Superfund Amendments and Reauthorization Act of 1986. The project was assigned under START Contract No. 68-S7-01-41, Task Order No. 0027.

This analysis was conducted using data collected from August 23, 2002 (Round 7) to December 22, 2003 (Round 15) from properties located within 0.25 mile of the smelter (house numbers 5, 6, 20, 22, and 24).

This analysis had the following objectives:

- (1) Evaluate differences in lead concentrations between front and backyards, and determine the most appropriate way to use data collected from individual quadrants in any statistical analysis of temporal trends.
- (2) Conduct a detailed trend analysis for both individual properties and all properties combined, and provide estimates of the potential rates of change in lead concentrations, as well as an assessment of uncertainties associated with these estimates.

Objective 1:

An assessment of lead concentrations in front versus backyards was conducted using qualitative (graphical) and quantitative (paired-difference statistical tests) methods. Figure 1 presents outlier box plots (also referred to as box-and-whisker plots) of lead concentrations by sampling round for individual properties and for all properties combined in the left panel of each plot. In these plots, lead measurements from front yards are shown as open circles, and concentrations in backyards are shown as solid circles. The right-hand panel of each plot shows front versus backyard comparisons for all rounds combined. Horizontal lines that transect each box are the arithmetic mean concentrations. The bottom and top of each box represent the 25th and 75th percentiles of the data, respectively. The area between the 25th and 75th percentiles is referred to as the interquartile range (IQR). The horizontal line drawn across the

interior of the boxes represents the median or 50th percentile. The upper and lower bounds of the whiskers represent the highest and lowest values, respectively, that are not considered statistical outliers. Points falling above the whiskers are considered “high outliers,” defined as values greater than the 75th percentile plus 1.5 times the IQR. “Low outliers” are defined as values less than the 25th percentile minus 1.5 times the IQR.

Statistical analysis of front and backyard concentrations was conducted using a paired-difference approach (Figure 2). This approach effectively evaluates the average difference in front versus backyard concentrations over all rounds of sampling. To conduct paired-difference testing, the average front and backyard concentrations (average of the two quadrants) are first calculated for each round of sampling. Next, the difference (front-backyard) between the two mean concentrations is calculated for each round. To conclude that no difference exists between front and backyard concentrations, the average difference over all sampling rounds should not differ statistically from zero, though some positive differences and some negative differences are expected. Figure 2 provides box plots, frequency histograms, and quantile tables of observed differences calculated for each sampling round. Figure 2 also provides results of both parametric (paired-difference t-test) and nonparametric (signed-rank test) tests of the two-sided null hypothesis (H_0) that the average difference is zero. If the probabilities associated with either of the two-sided tests (that is, “Prob > |t|” under the “Test Statistic” column in the embedded table) are less than or equal to 0.05 (5 percent), H_0 is rejected with the conclusion that lead concentrations are statistically different between front and backyards.

Conclusions for Objective 1. The statistical tests show that front yard concentrations exceed backyard concentrations for property numbers 22 and 24, and that backyard concentrations exceed front yard concentrations for property number 5. No statistically significant difference in concentrations was shown for property numbers 6 and 20. Examination of Figure 1 and results of the paired-difference tests suggest no consistent spatial pattern among properties when data are divided into front and backyard groupings. This is not surprising, as these groupings are not based on absolute orientations of front or backyards (or in effect, the relative degree of exposure of each area to air emissions of lead). A more technically defensible way to analyze the data at the scale of an individual quadrant (or pair of quadrants) would be to create groupings that reflect the relative degree of exposure to lead emissions. Under this approach, a front or backyard for an individual property would be grouped, as “exposed” or “relatively unexposed,” and statistical analysis would proceed accordingly.

Objective 2:

Previous trend analyses conducted using a nonparametric test for monotonic trends (Mann-Kendall test) revealed significant increasing trends in lead concentrations for house numbers 5, 20, and 22 for samples collected from Round 7 through 14. To estimate potential rates of increase in lead concentrations, this analysis was repeated using linear regression analysis and additional data collected during Round 15. For this analysis, the earliest sampling date for Round 7 (August 23, 2002) was set at day = 0, and all subsequent sampling dates were converted to the number of days from this initial sampling date.

Figure 3 presents the results of regression analysis conducted for individual properties using two approaches. The first approach treats the data from individual quadrants as independent measurements. These results are presented as the set of plots at the left in Figure 3. Open and solid circles are used to represent data collected from front and backyards, respectively. The second approach calculates a median concentration for all four quadrants and regresses this against the time variable. As shown by the previous tests for monotonic trends, significant increasing trends appear for house numbers 5, 20, and 22. The trend results are significant for both the case in which individual measurements were used for each quadrant and the case in which only the median concentrations were used in the regression.

Figure 4 presents the results of regression analysis with all properties within 0.25 mile of the smelter combined. Three approaches are used in the analysis for the combined data:

- (1) All measurements for all quadrants and houses are used in the regression of time (days) versus concentration.
- (2) The medians for each round are calculated for individual houses, and the medians for all five houses are used in the regression of time (days) versus concentration.
- (3) The grand median of all houses combined is calculated for each round, and the grand median for each round is used in the regression of time (days) versus concentration.

A significant increasing trend results from application of all three approaches using the data for all houses combined.

Uncertainty analysis. Graphical approaches as well as parametric and nonparametric statistical trend tests provide clear evidence of an increasing trend in lead concentrations for selected properties, and for the scenario in which all properties are treated as a single area. Therefore, estimating the potential rate at which lead concentrations may be increasing is of interest.

Estimating potential rates of increase in lead concentrations is not straightforward, and a number of uncertainties are associated with the types of trend analyses routinely conducted using environmental data sets. These uncertainties include:

- 1) Spatial variability in measured lead concentrations. That is, rates of increase are unlikely to be uniform over large areas, and different ways of grouping properties in this type of analysis are likely to yield different estimates. Moreover, different ways are available to treat replicate data for individual quadrants for each property. Figures 3 and 4 show clear differences in the relative spread of the data (that is, variability or noise) around the fitted regression line in cases where data for individual quadrants are included as independent measurements as opposed to representation by a measure of central tendency (median).
- 2) Selection of an appropriate statistical model for quantifying rates of increase. Linear models are the simplest to evaluate, but more complex non-linear models may better represent the true relationship between concentration and time.
- 3) Imprecision in measured concentrations associated with use of XRF analysis.

To better quantify the uncertainty associated with estimates of the rate of increase provided in Figures 3 and 4, two-sided 95-percent confidence limits were calculated for individual slope factors (regression coefficients) under several scenarios (Table 1). Two scenarios were considered. The first was a “worst case” scenario in which confidence limits were calculated for the property showing the most rapid increase in lead concentrations (house 20). The second was the scenario in which all properties were combined and treated as a single area. Regression coefficients (slope factors) in Table 1 are reported in units of milligrams per kilogram (mg/kg) per day as well as mg/kg per month. Confidence limits in Table 1 are expressed as monthly rates of increase.

It should be noted that EPA QA/G-9 (EPA 2000) provides several cautionary statements about applying regression analysis to environmental data. EPA (2000) notes that regression analysis is sensitive to extreme values (outliers), and that it is not well suited for handling censored data (that is, data below the detection limit). The validity of regression analysis also depends on two key assumptions: normally distributed errors and constant variance. EPA (2000) states that verifying these assumptions may be difficult or burdensome in routine practice, and suggests that regression analysis may be most useful as an “informal, quick, and easy screening tool for identifying strong linear trends.” In the present analysis, the distribution of errors (that is, the deviation of individual measurements relative to the fitted regression line) was evaluated using residual plots. A plot of the residuals, or deviation of individual measurements from the “line of best fit,” allows a rapid qualitative check of the assumption of normally distributed

errors. Residual plots appear at the bottom of the individual regression plots provided in Figures 3 and 4. Straightforward visual inspection of the residual plots does not provide evidence of serious departure from the assumption of normally distributed errors.

EPA (2000) also recommends a nonparametric approach for estimating slopes (Sen's Slope Estimate) as an alternative to linear regression. Calculation of Sen's Slope Estimate by hand is very tedious and not provided by most commercial statistical software packages. Sen's Slope Estimate was calculated in this analysis using a customized computer model similar to that used for running the Mann-Kendall test for monotonic trends. Calculations were performed for the worst-case scenario and the scenario in which the data for all houses are pooled. For both scenarios, only the data for the medians of all quadrants were used. These calculations yielded point estimates for the slope of 22 and 8 mg/kg per month for the worst-case and pooled scenarios, respectively.

Conclusions for Objective 2. The rates of increase shown in Table 1 should be viewed as rough estimates of the true rates of increase, given the uncertainties previously discussed. While regression analysis typically should not be used to predict future observations (that is, predictions should be limited to the range of the observed data), a coarse approximation of lead deposition over time could be made by assuming that the rate of increase predicted using the existing data would remain relatively constant. If so, the confidence limits calculated for the worst case scenario (increases as high as 6 to 11 mg/kg per month) and the scenario with all properties combined (increases as high as 1 to 6 mg/kg per month) could be used to estimate the time until average concentrations of lead for individual properties (or area-wide averages) will exceed the 400 mg/kg preliminary remediation goal (PRG) for lead. Collecting additional data over time and/or increasing sampling densities for individual rounds would increase the precision and accuracy of these estimates.

REFERENCE

U.S. Environmental Protection Agency (EPA). 2000. "Guidance for Data Quality Assessment: Practical Methods for Data Analysis: EPA QA/G-9, QA97 Version." EPA/600/R-96/084. Office of Research and Development. Washington, DC. July.

TABLE 1
RESULTS OF LINEAR REGRESSION ANALYSIS FOR SELECTED SCENARIOS

Scenario	Sample Size	Sampling Interval		Regression Coefficients for Days Versus Concentration			95 Percent Confidence Limits for Monthly Increase in Lead Concentrations	
		Initial Date	Final Date	Intercept	Slope ¹	S.E. (Slope)	Lower C.L.	Upper C.L.
Worst Case- House No. 20 (all quadrants treated as independent measurements)	32	8/26/2002	12/22/2003	72.74	0.28 (8.51)	0.04	6.12	10.80
Worst Case- House No. 20 (median of all quadrants used for each sampling round)	8	8/26/2002	12/22/2003	63.26	0.32 (9.53)	0.03	7.26	11.36
All Houses Combined (all quadrants treated as independent measurements)	148	8/23/2002	12/22/2003	85.98	0.11 (3.26)	0.03	1.70	4.81
All Houses Combined (medians for each house used for each round)	37	8/23/2002	12/22/2003	77.39	0.11 (3.43)	0.04	1.09	5.69
All Houses Combined (medians of all houses and quadrants used for each round)	9	8/23/2002	12/22/2003	68.74	0.12 (3.65)	0.02	2.03	5.00

Notes:

C.L. Multiply the slope by 30 to get a point estimate of the monthly increase
Confidence Limits
S.E. Standard error of estimate
1 Numbers in parentheses are slope factors expressed as mg/kg per month

FIGURE 1

BOX PLOTS OF LEAD CONCENTRATIONS BY ROUND AND COMPARISON OF FRONT AND BACK YARDS FOR ALL ROUNDS COMBINED

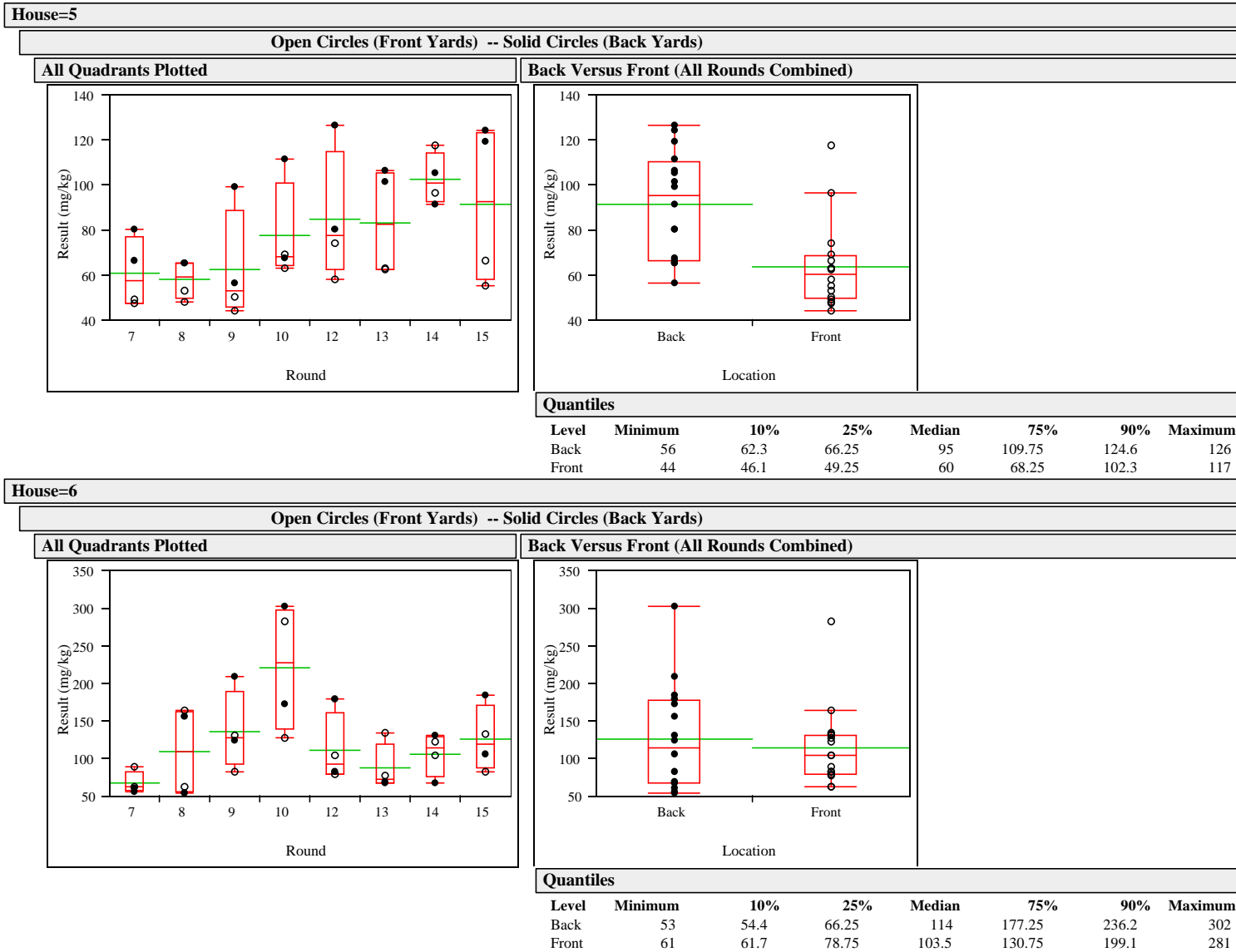


FIGURE 1 (CONTINUED)

BOX PLOTS OF LEAD CONCENTRATIONS BY ROUND AND COMPARISON OF FRONT AND BACK YARDS FOR ALL ROUNDS COMBINED

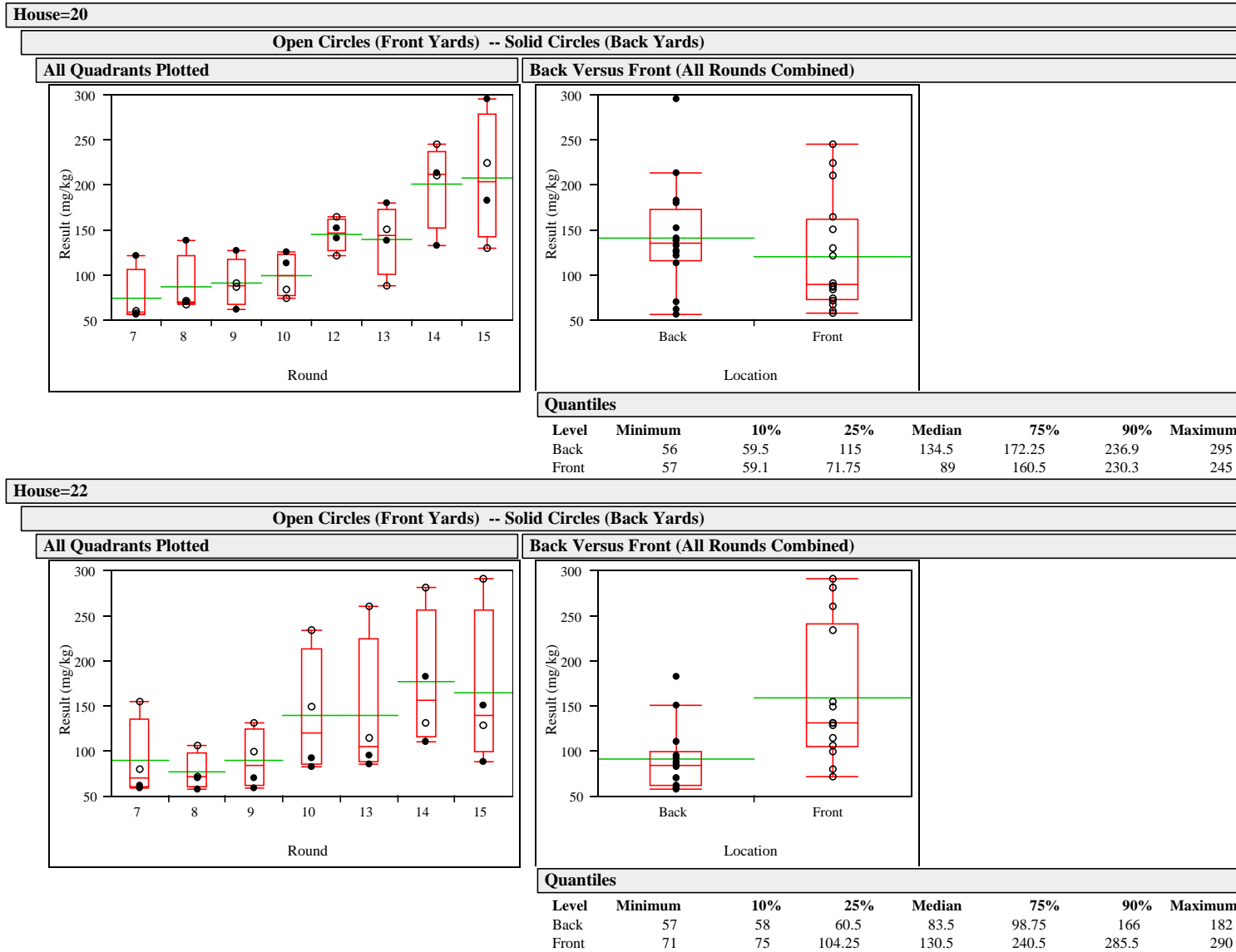


FIGURE 1 (CONTINUED)

BOX PLOTS OF LEAD CONCENTRATIONS BY ROUND AND COMPARISON OF FRONT AND BACK YARDS FOR ALL ROUNDS COMBINED

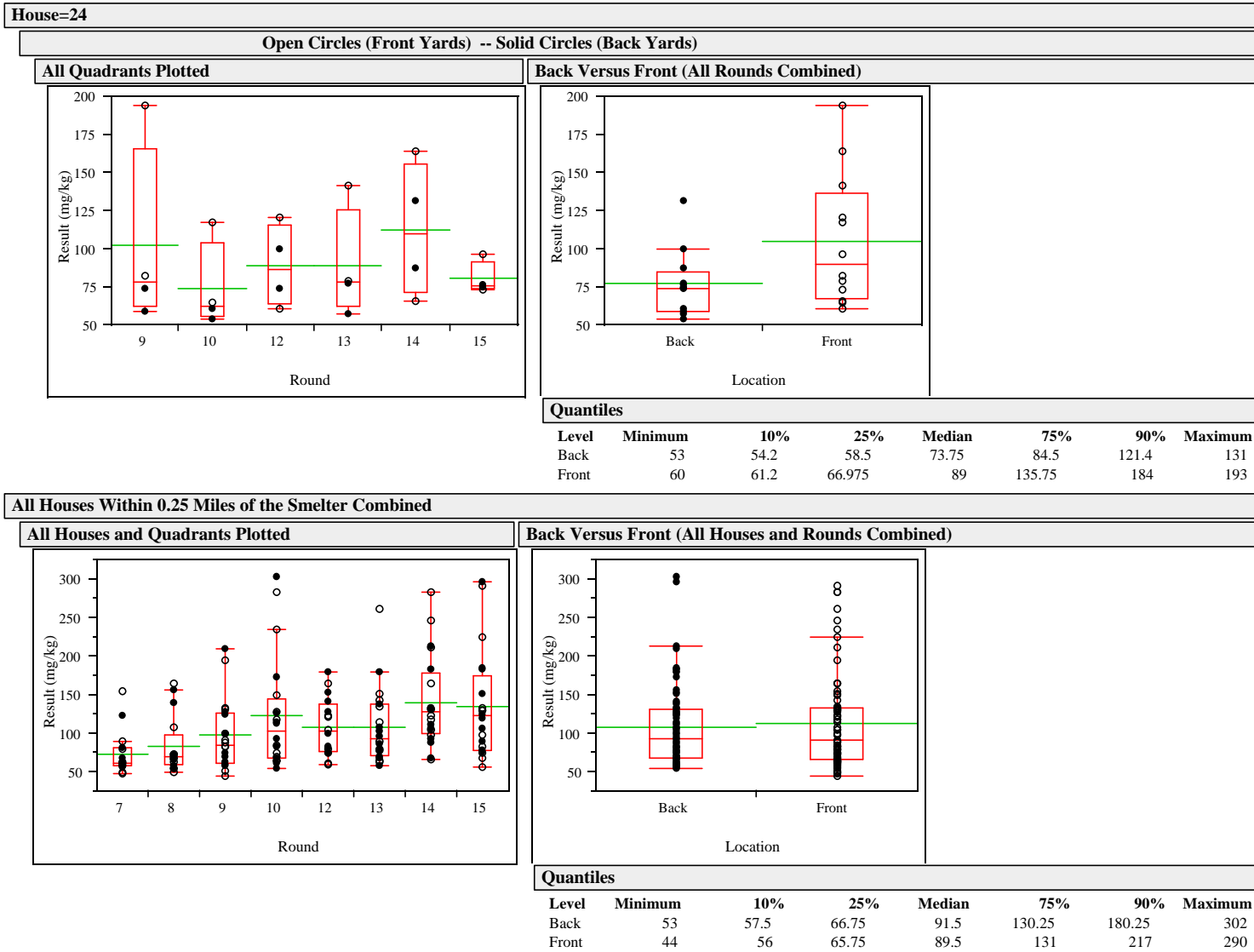


FIGURE 2

RESULTS OF PAIRED-DIFFERENCE TESTS COMPARING LEAD CONCENTRATIONS IN FRONT AND BACK YARDS

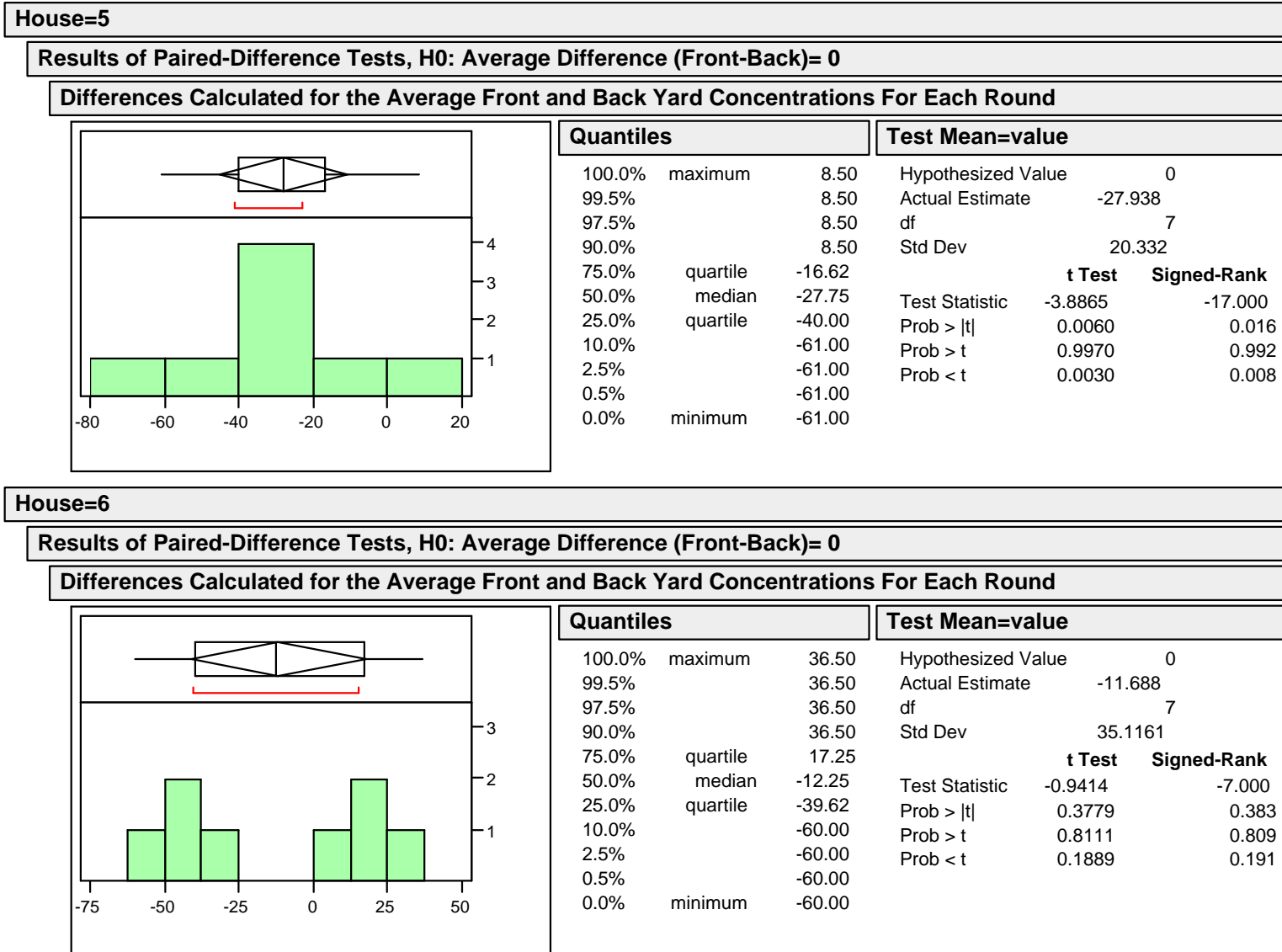


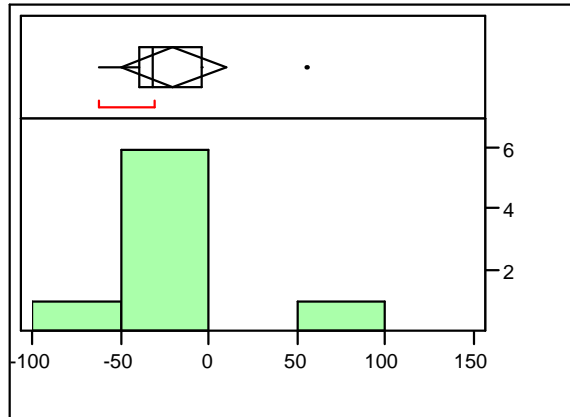
FIGURE 2 (CONTINUED)

RESULTS OF PAIRED-DIFFERENCE TESTS COMPARING LEAD CONCENTRATIONS IN FRONT AND BACK YARDS

House=20

Results of Paired-Difference Tests, H0: Average Difference (Front-Back)= 0

Differences Calculated for the Average Front and Back Yard Concentrations For Each Round

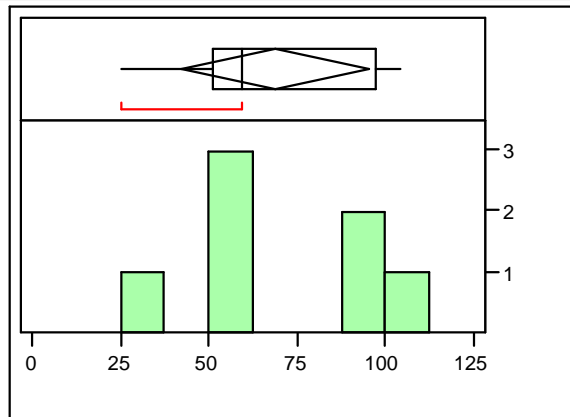


Quantiles			Test Mean=value	
100.0%	maximum	55.50	Hypothesized Value	0
99.5%		55.50	Actual Estimate	-19.875
97.5%		55.50	df	7
90.0%		55.50	Std Dev	35.8924
75.0%	quartile	-4.00	t Test	Signed-Rank
50.0%	median	-32.25	Test Statistic	-1.5662
25.0%	quartile	-39.75	Prob > t	0.1613
10.0%		-62.00	Prob > t	0.9194
2.5%		-62.00	Prob < t	0.0806
0.5%		-62.00		
0.0%	minimum	-62.00		

House=22

Results of Paired-Difference Tests, H0: Average Difference (Front-Back)= 0

Differences Calculated for the Average Front and Back Yard Concentrations For Each Round



Quantiles			Test Mean=value	
100.0%	maximum	104.00	Hypothesized Value	0
99.5%		104.00	Actual Estimate	69
97.5%		104.00	df	6
90.0%		104.00	Std Dev	28.7417
75.0%	quartile	97.00	t Test	Signed-Rank
50.0%	median	59.50	Test Statistic	6.3516
25.0%	quartile	51.00	Prob > t	0.0007
10.0%		25.00	Prob > t	0.0004
2.5%		25.00	Prob < t	0.9996
0.5%		25.00		
0.0%	minimum	25.00		

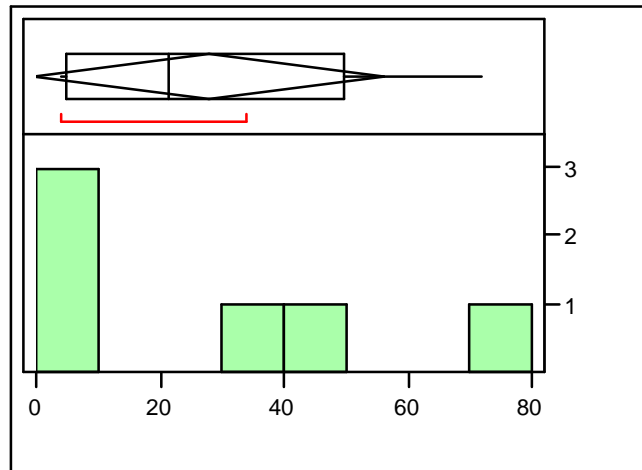
FIGURE 2 (CONTINUED)

RESULTS OF PAIRED-DIFFERENCE TESTS COMPARING LEAD CONCENTRATIONS IN FRONT AND BACK YARDS

House=24

Results of Paired-Difference Tests, H0: Average Difference (Front-Back)= 0

Differences Calculated for the Average Front and Back Yard Concentrations For Each Round



Quantiles			Test Mean=value		
100.0%	maximum	72.000	Hypothesized Value	0	
99.5%		72.000	Actual Estimate	27.7667	
97.5%		72.000	df	5	
90.0%		72.000	Std Dev	26.9968	
75.0%	quartile	49.875		t Test	Signed-Rank
50.0%	median	21.550	Test Statistic	2.5193	10.500
25.0%	quartile	4.750	Prob > t	0.0532	0.031
10.0%		4.000	Prob > t	0.0266	0.016
2.5%		4.000	Prob < t	0.9734	0.984
0.5%		4.000			
0.0%	minimum	4.000			

FIGURE 3

LINEAR REGRESSION ANALYSIS CONDUCTED USING DATA FOR INDIVIDUAL QUADRANTS AND FOR THE MEDIAN OF ALL QUADRANTS

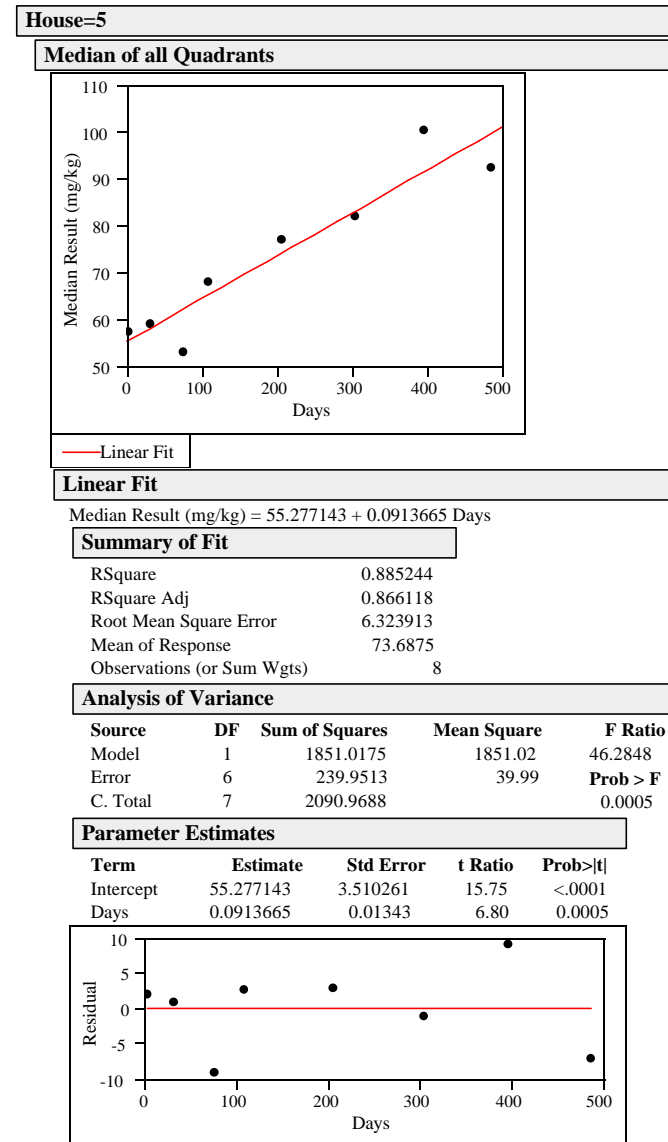
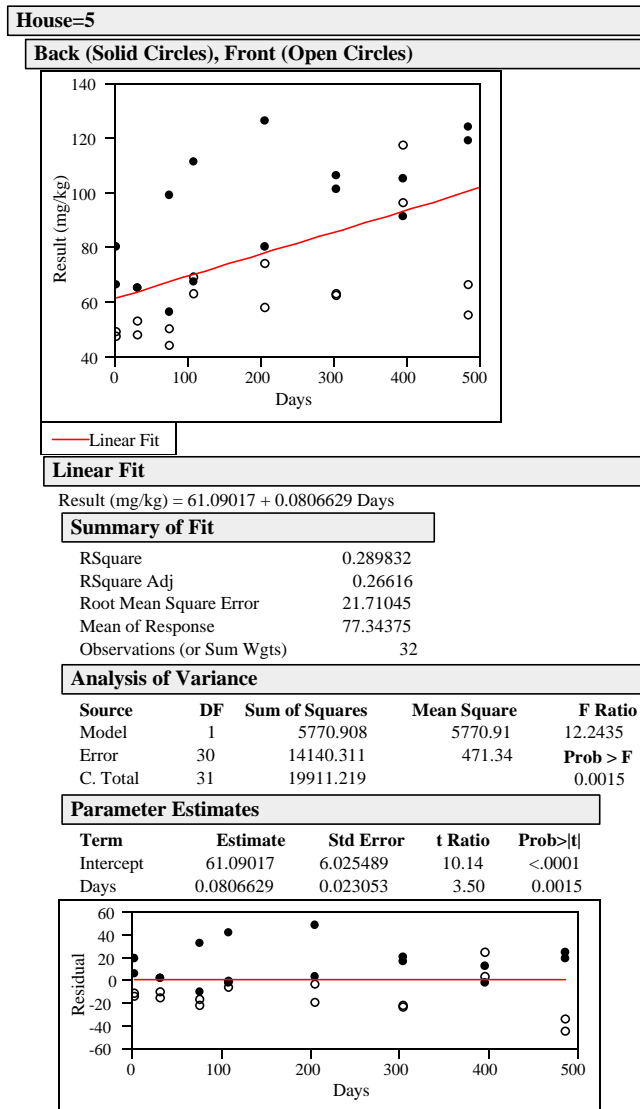


FIGURE 3 (CONTINUED)

LINEAR REGRESSION ANALYSIS CONDUCTED USING DATA FOR INDIVIDUAL QUADRANTS AND FOR THE MEDIAN OF ALL QUADRANTS

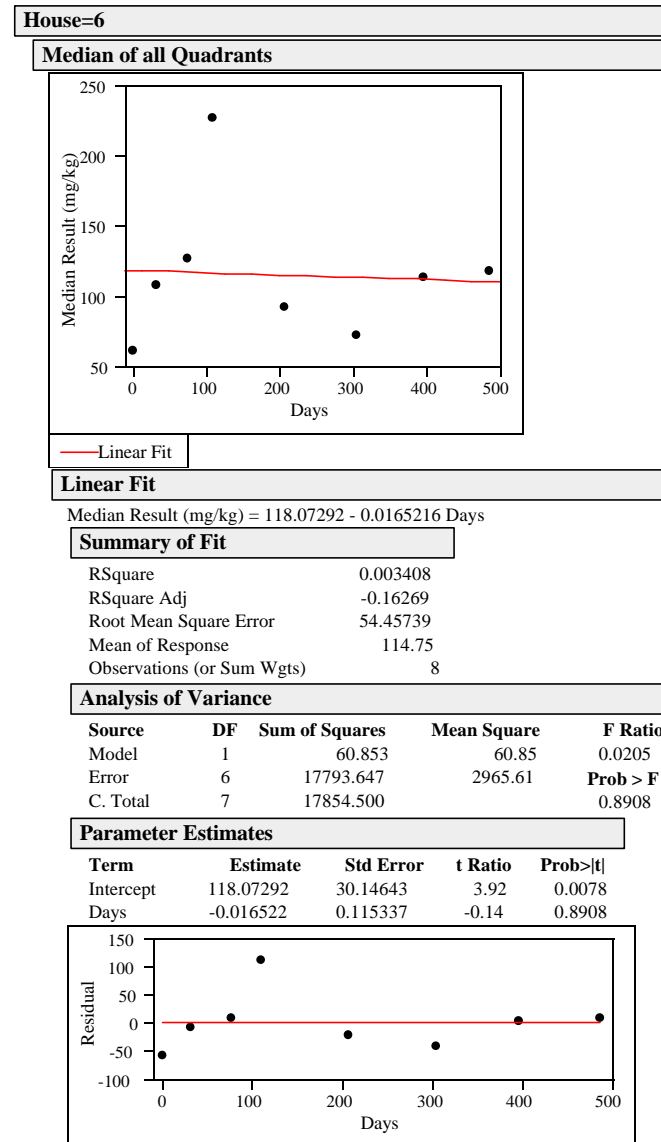
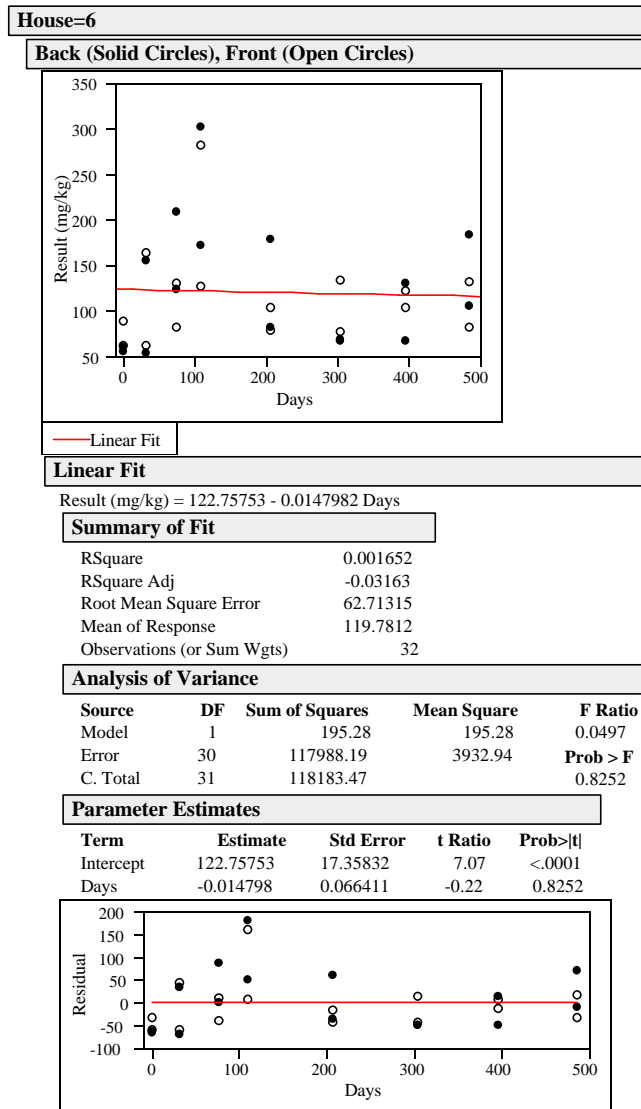


FIGURE 3 (CONTINUED)

LINEAR REGRESSION ANALYSIS CONDUCTED USING DATA FOR INDIVIDUAL QUADRANTS AND FOR THE MEDIAN OF ALL QUADRANTS

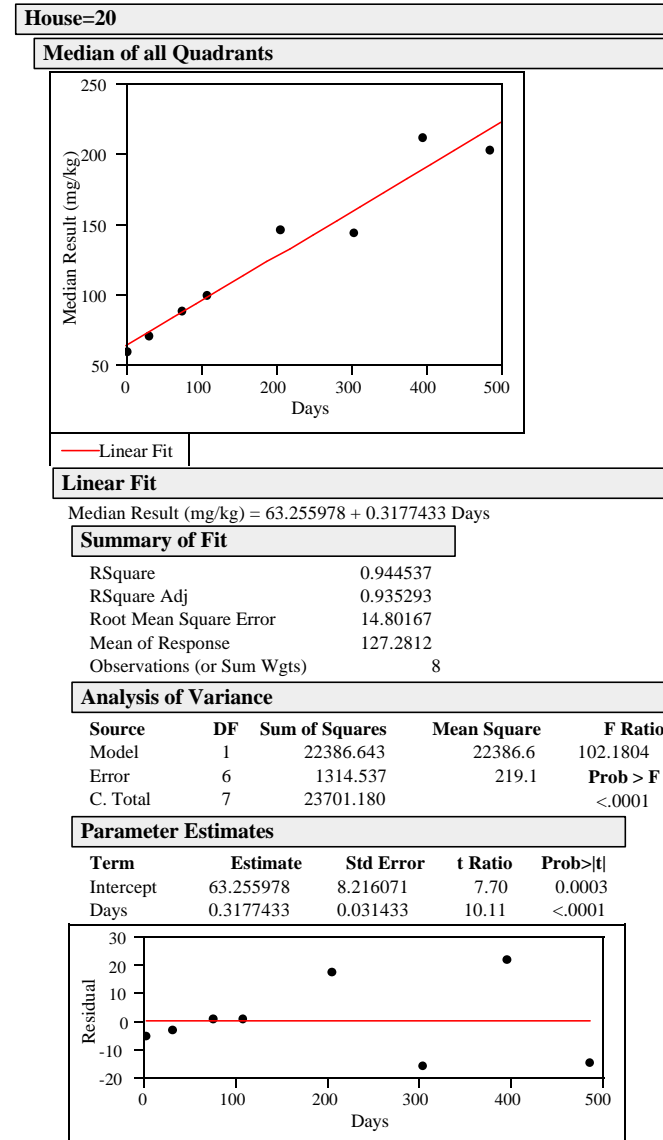
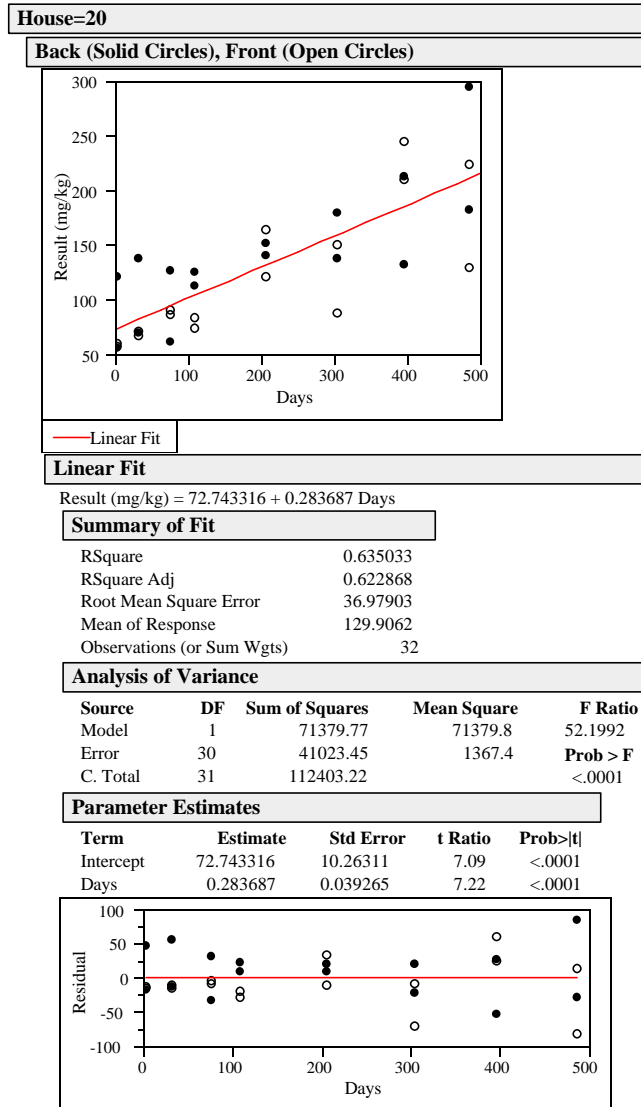


FIGURE 3 (CONTINUED)

LINEAR REGRESSION ANALYSIS CONDUCTED USING DATA FOR INDIVIDUAL QUADRANTS AND FOR THE MEDIAN OF ALL QUADRANTS

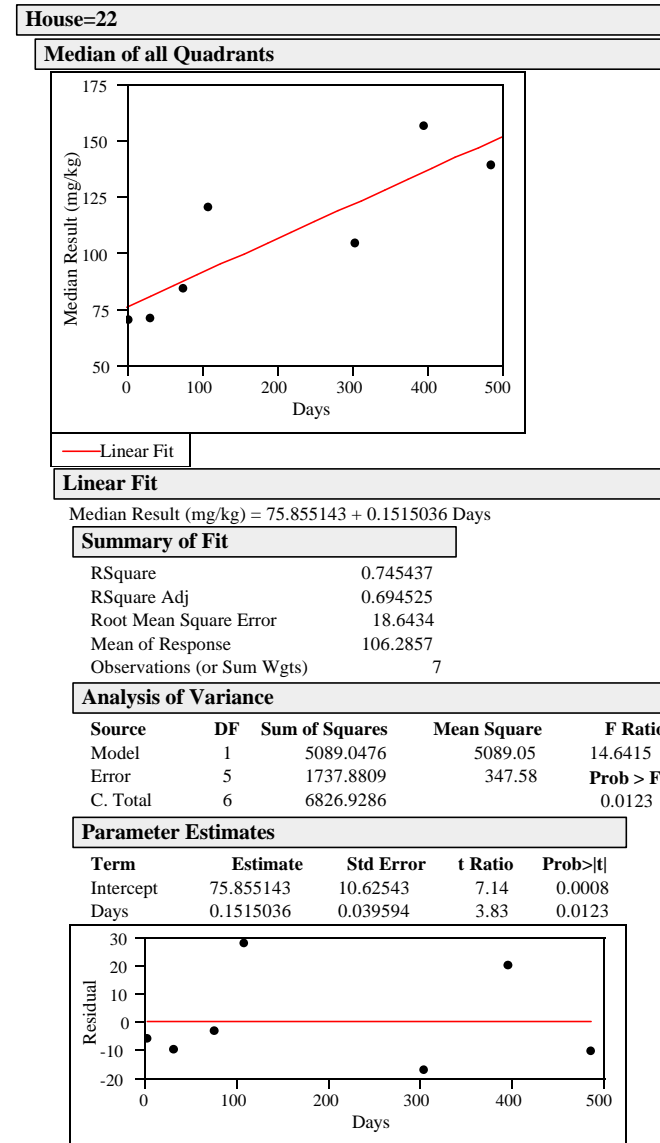
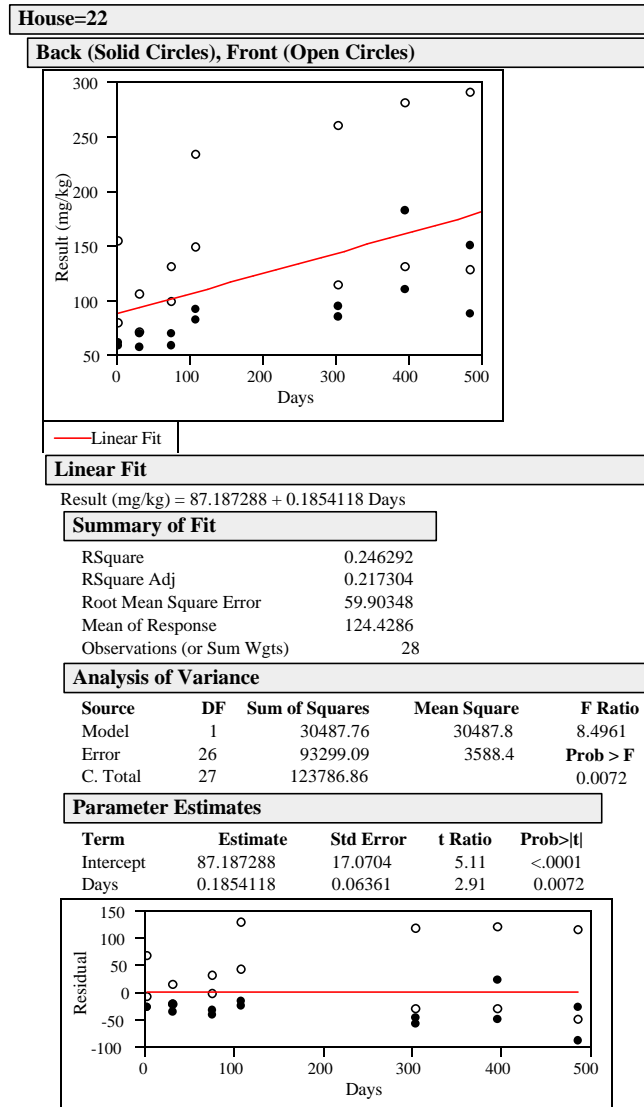


FIGURE 3 (CONTINUED)

LINEAR REGRESSION ANALYSIS CONDUCTED USING DATA FOR INDIVIDUAL QUADRANTS AND FOR THE MEDIAN OF ALL QUADRANTS

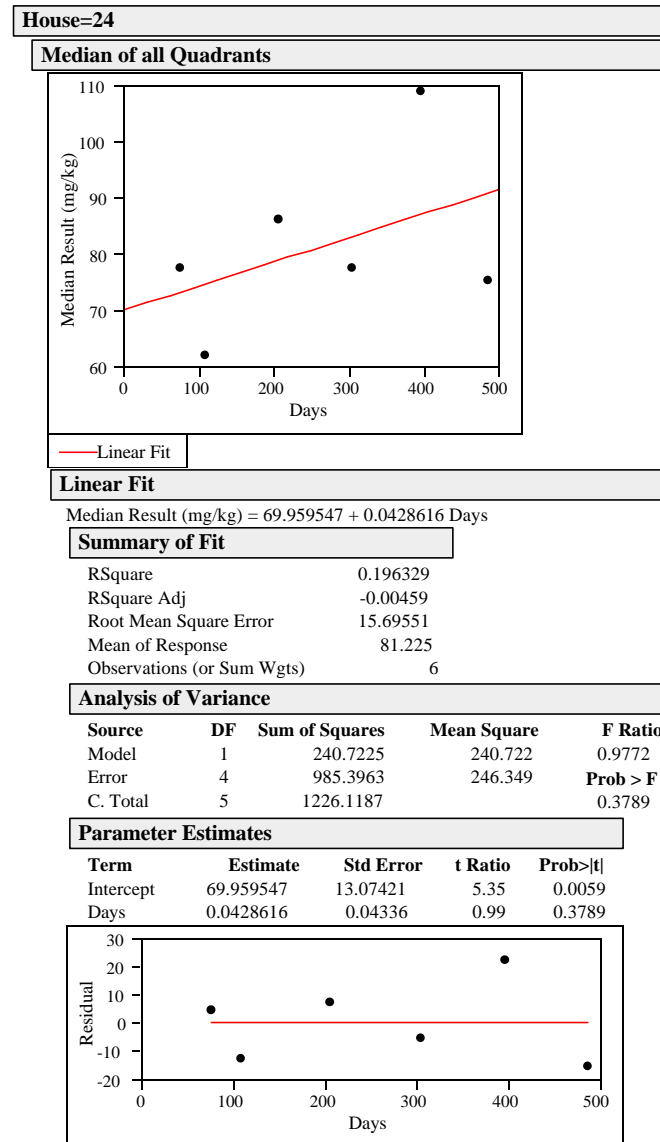
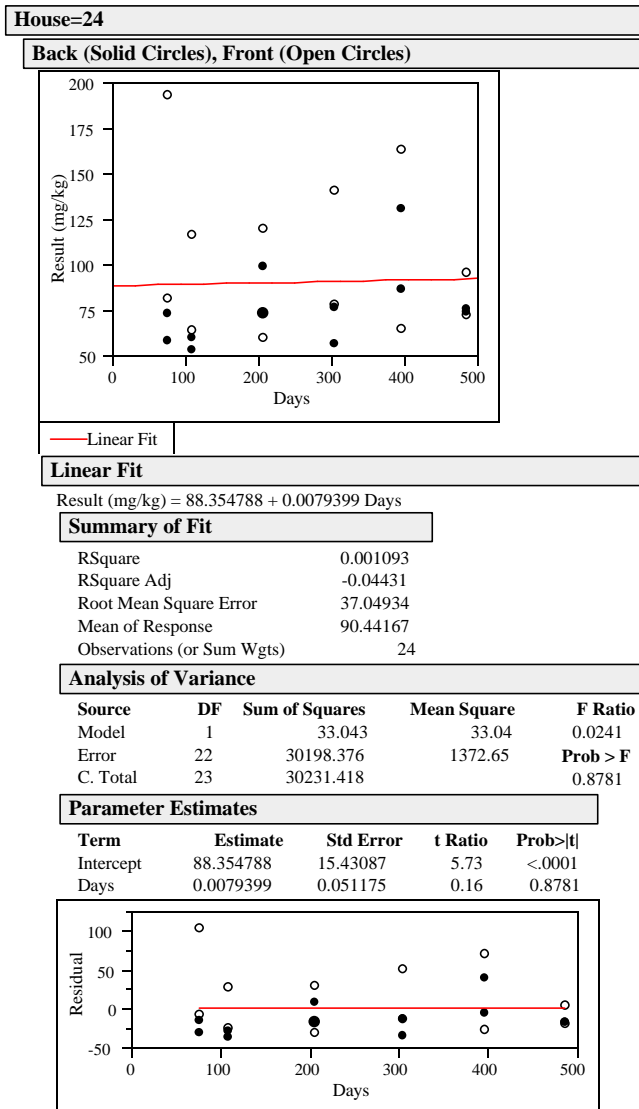
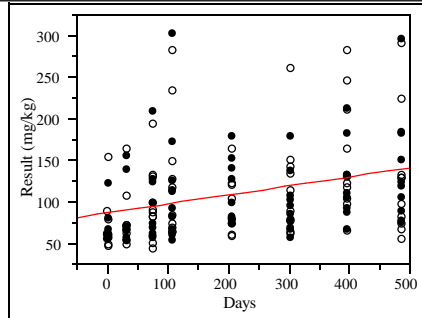


FIGURE 4

REGRESSION ANALYSIS FOR ALL HOUSES POOLED USING DATA FOR INDIVIDUAL QUADRANTS, THE MEDIAN OF ALL QUADRANTS FOR EACH HOUSE, AND THE OVERALL MEDIAN OF ALL QUADRANTS AND HOUSES COMBINED

All Quadrants and Houses Combined



Linear Fit

Linear Fit

Result (mg/kg) = 85.977667 + 0.1085883 Days

Summary of Fit

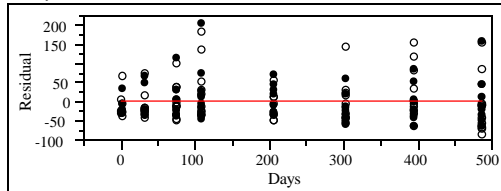
RSquare	0.104361
RSquare Adj	0.098226
Root Mean Square Error	53.65106
Mean of Response	108.9162
Observations (or Sum Wgts)	148

Analysis of Variance

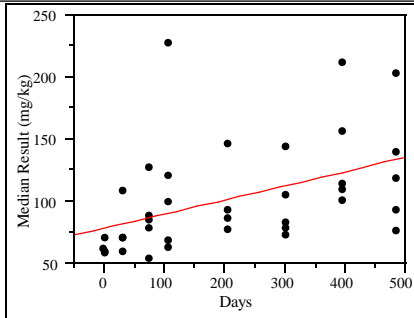
Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	48968.12	48968.1	17.0121
Error	146	420251.64	2878.4	Prob > F
C. Total	147	469219.76		<.0001

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	85.977667	7.097784	12.11	<.0001
Days	0.1085883	0.026327	4.12	<.0001



Medians Plotted for Individual Houses for Each Date



Linear Fit

Linear Fit

Median Result (mg/kg) = 77.391699 + 0.1143305 Days

Summary of Fit

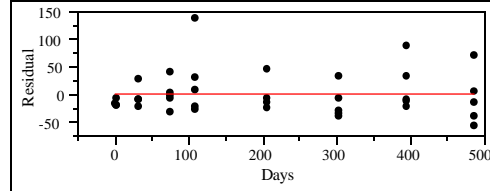
RSquare	0.201839
RSquare Adj	0.179035
Root Mean Square Error	39.1574
Mean of Response	101.5432
Observations (or Sum Wgts)	37

Analysis of Variance

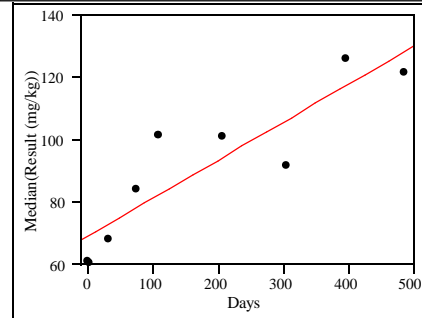
Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	13570.985	13571.0	8.8508
Error	35	53665.581	1533.3	Prob > F
C. Total	36	67236.566		0.0053

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	77.391699	10.36068	7.47	<.0001
Days	0.1143305	0.03843	2.98	0.0053



Grand Median of all Houses and Quadrants Plotted



Linear Fit

Linear Fit

Median(Result (mg/kg)) = 68.744595 + 0.1217734 Days

Summary of Fit

RSquare	0.801147
RSquare Adj	0.772739
Root Mean Square Error	11.64334
Mean of Response	90.55556
Observations (or Sum Wgts)	9

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	3823.2503	3823.25	28.2018
Error	7	948.9719	135.57	Prob > F
C. Total	8	4772.2222		0.0011

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	68.744595	5.650786	12.17	<.0001
Days	0.1217734	0.022931	5.31	0.0011

